

Minimal-variance distributed scheduling under strict demands and deadlines

Yorie Nakahira
Caltech
ynakahir@caltech.edu

Andres Ferragut
Universidad ORT Uruguay
ferragut@ort.edu.uy

Adam Wierman
Caltech
adamw@caltech.edu

ABSTRACT

Many modern schedulers can dynamically adjust their service capacity to match the incoming workload. At the same time, however, variability in service capacity often incurs operational and infrastructure costs. In this abstract, we characterize an optimal distributed algorithm that minimizes service capacity variability when scheduling jobs with deadlines. Specifically, we show that *Exact Scheduling* minimizes service capacity variance subject to strict demand and deadline requirements under stationary Poisson arrivals. Moreover, we show how close the performance of the optimal distributed algorithm is to that of the optimal centralized algorithm by deriving a competitive-ratio-like bound.

Keywords

Deadline scheduling, Service capacity control, Exact Scheduling, Online distributed algorithm

1. INTRODUCTION

Traditionally, the scheduling literature has assumed a static (fixed) service capacity. However, it is increasingly common for modern applications to have the ability to dynamically adjust their service capacity in order to match the current demand. For example, when using cloud computing services, one can modify the total computing capacity by changing the number of computing instances and their speeds. Power distribution networks can also adapt the energy supply to match the energy demand as it changes over time.

The ability to adapt service capacity dynamically gives rise to challenging new design questions. In particular, how to maintain *predictability* and *stability* of service capacity is of great importance in such applications since peaks and fluctuations often come with significant costs [1]. This trend is especially true for the examples of cloud computing and power distribution networks mentioned above. Cloud content providers prefer stable and predictable service capacity because on-demand contracts for compute instances (e.g., Amazon EC2 and Microsoft Azure) are typically more expensive than long-term contracts. Additionally, large fluctuations in service capacity induce unnecessary power consumption and infrastructure strain for computing equipment. The emerging load from electric vehicle charging stations also leads to similar challenges in power distribution networks. Charging stations require stability in power con-

sumption because fluctuations and large peaks in power use may strain the grid infrastructure and result in a high peak charge for the station operators. The stations also prefer predictable power consumption because purchasing power in real time is typically more expensive than purchasing in advance.

Thus, in situations where service capacity can be dynamically adjusted, an important design goal is to minimize the costs associated with variability in the service capacity while maintaining high quality of service. In this paper, we study this problem by minimizing the *variance* of the service capacity in systems where jobs arrive with demand and deadline requests. Our focus on service capacity variance is motivated by applications such as cloud computing and power distribution networks, where contracts often explicitly depend on service capacity variability, e.g., if a charging station participates in the regulation market, then costs/payments rely explicitly on the variance of the total capacity [2].

Although the literature on deadline scheduling is large and varied (see [3, 4] and references therein), the optimal algorithms are only known for certain niche cases such as deterministic worst-case settings [5], single server systems [6], and heavy traffic settings [7]. In particular, the problem of designing an *optimal* algorithm that minimizes service capacity variability while satisfying service quality constraints, i.e., meeting demands and deadlines, has remained open. Solving this problem is a challenging task due to the heterogeneous constraints (diversity in service requests) and the size of the state and decision space (the number of possible remaining job profile configurations and feasible control policies).

The goal of this work is to characterize the distributed scheduling algorithm that minimizes the variance of service capacity subject to service quality constraints, e.g., meeting job deadlines and satisfying job demands. Our focus is on *distributed* algorithms since implementing centralized algorithms is likely to be prohibitively slow and costly in large-scale service systems today. From cloud computing to power distribution networks, such systems are unlikely to be able to access global information about every job and server in the system when deciding the service rate of each job/server. Therefore, distributed algorithms are a necessity to enable large-scale implementation.

Our Contributions. In this work, we characterize the optimal distributed policy by using tools from optimization and control theory. Specifically, we show that *Exact Scheduling* is the optimal distributed algorithm, i.e., it minimizes

the stationary variance of the service capacity among all distributed policies that strictly satisfy the service requirements. Exact Scheduling is a classical algorithm that works by finishing job service *exactly* at their deadlines using a constant service rate [1, 3, 8]. Given that our results focus on distributed algorithms, we also study how the optimal distributed algorithm performs compared with the optimal centralized algorithm, which may provide better performance in theory but requires prohibitively expensive computation to find in practice. To answer this question, we derive a closed-form bound on the performance degradation due to using a distributed algorithm. The bound suggests that, when the sojourn time is a deterministic variable, Exact Scheduling attains the optimal trade-off between service capacity variance and total remaining demand variance achievable by any centralized algorithms.

2. PROBLEM FORMULATION

We consider a setting in which a service system dynamically adjusts its capacity to serve jobs that arrive randomly with heterogeneous service requirements. We use a continuous time model where $t \in \mathbb{R}_+$ denotes a point in time. Each job, indexed by $k \in \mathcal{V} = \{1, 2, \dots\}$, is characterized by a random arrival time a_k , a random service demand σ_k , and a random sojourn time $\tau_k (\geq \sigma_k)$.¹ To formalize the scheduler design procedure into an optimization problem, we introduce below the arrival profiles, the service profiles, the system dynamics, and the design objectives in detail.

Arrival profiles. We represent the set of jobs as a marked point process $\{(a_k; \sigma_k, \tau_k)\}_{k \in \mathcal{V}}$ in $\mathbb{R}_+ \times S$, where the arrival times $a_k \in \mathbb{R}_+$ are the set of points, and the service requirements $(\sigma_k, \tau_k) \in S$ are the set of marks. We assume that the marked point process is a stationary independently marked Poisson Point Process, which is defined by a locally finite non-null intensity measure $\Lambda(a)$ on \mathbb{R}_+ and a mark density measure $f(\sigma, \tau)$ on S [9].² This also implies that $\{(a_k; \sigma_k, \tau_k)\}_{k \in \mathcal{V}}$ is a Poisson Point Process on $\mathbb{R}_+ \times S$ with an intensity measure $f(\sigma, \tau)\Lambda$. Intuitively, $\int_A f(\sigma, \tau)\Lambda d\sigma d\tau$ is the average rate at which jobs with service requirement $(\sigma, \tau) \in A \subset S$ arrive. We additionally assume that S is bounded, and $S \subset \{(\sigma, \tau) : \tau \geq \sigma \geq 0\}$.¹

Service profiles. The service system works on each job $k \in \mathcal{V}$ with a *service rate* $r_k(t) \geq 0$. The service rate can take non-zero values only when the job sojourns in the system, i.e., $r_k(t) = 0$ for any $t \notin [a_k, a_k + \tau_k]$. To meet the service demand of job k , its service rate must satisfy

$$\int_{a_k}^{a_k + \tau_k} r_k(t) dt = \sigma_k, \quad k \in \mathcal{V}. \quad (1)$$

Without loss of generality, $r_k(t) = 1$ is assumed to be the maximum rate: that is, $r_k(t)$ can take any values in $[0, 1]$, and $r_k < 1$ corresponds to throttling down service speed at the expense of prolonging job completion times. The above sojourn time and maximum rate constraints can be jointly written as

$$0 \leq r_k(t) \leq \mathbf{1}\{t \in [a_k, a_k + \tau_k]\}, \quad (2)$$

¹The condition $\tau_k \geq \sigma_k$ requires a job to have a service demand that is no more than the maximum service (which is assume to be 1) times the sojourn time.

²Here, we use $(a; \sigma, \tau)$ to denote random variables and $(a_k; \sigma_k, \tau_k)$ to denote one realization of them in job k .

where $\mathbf{1}\{A\}$ denotes the indicator function for an event A . The resource consumption is defined by the (active) service capacity

$$P(t) = \sum_{k \in \mathcal{V}} r_k(t).$$

System dynamics. At each time $t \in \mathbb{R}_+$, job k has a remaining demand $x_k(t) = \sigma_k - \int_{a_k}^t r_k(h) dh$ and a remaining time $y_k(t) = a_k + \tau_k - t$. The set of remaining jobs in the system can be considered as a Point Process $\{(x_k(t), y_k(t))\}_k$ in \mathbb{R}^2 , where x -axis is the remaining demand and y -axis is the remaining time (see Fig 1). At time t , each point(job) moves with velocity $-r_k(t)$ in x -axis and velocity -1 in y -axis.

Scheduling algorithms. An online scheduling algorithm decides the service rates in real-time without using the future job arrival information. For scalability, we additionally restrict our attention to the following form of *distributed* algorithms which decide the service rate of a job only using its own information:

$$r_k(t) = u(x_k(t), y_k(t)) \geq 0. \quad (3)$$

Here, $u : \mathbb{R}_+^2 \rightarrow \mathbb{R}_+$ is a deterministic function of the remaining demand $x_k(t)$ and the remaining time $y_k(t)$ of each job k at time t . Under any policy of the form (3), the set of jobs remaining in the system converges to a stationary distribution. This stationary distribution is a Spatial Poisson Point Process with intensity measure $\lambda(x, y)$ satisfying

$$0 = \frac{\partial}{\partial x}(\lambda(x, y)u(x, y)) + \frac{\partial}{\partial y}\lambda(x, y) + \Lambda f(x, y), \quad (4)$$

where x is the remaining demand and y is the remaining time. Because the remaining job distribution converges to a stationary distribution, $P(t)$ also converges to a stationary distribution.³

Design objectives. We consider minimizing service capacity variability under hard demand and deadline constraints:³

$$\underset{u: (1)(2)(3)(4)}{\text{minimize}} \quad \text{Var}(P), \quad (5)$$

where $\text{Var}(P)$ is a functional of u and $\lambda(\sigma, \tau)$ satisfying (4). The optimization problem (5) has demand constraints as in (1) and service rate constraints as in (2), and the optimization variable u is constrained by (3) to be a distributed algorithm.

3. OPTIMAL SCHEDULING

When the goal is to minimize $\text{Var}(P)$, it is worth noting that peaks in service rate amplifies the uncertainties in the future arrivals, which in turn produce large variance in $P(t) = \sum_{k \in \mathcal{V}} r_k(t) = \sum_{k \in \mathcal{V}} u(x_k(t), y_k(t))$. This observation motivate us to use a flat service rate (Fig 1), which is achieved by the scheduling policy

$$u(x, y) = \begin{cases} \frac{x}{y}, & \text{if } y > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

This policy is known as Exact Scheduling and works by finishing all jobs *exactly* at their deadlines using a constant

³We use $\mathbb{E}[P]$ and $\text{Var}(P)$ to represent the stationary mean and variance of a stochastic process $P(t)$.

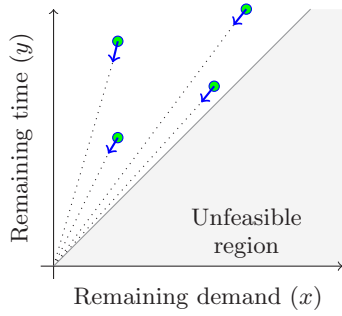


Figure 1: *Exact scheduling depicted in the space of remaining demand x and remaining time y .*

service rate. It is also highly scalable because it is distributed and asynchronous, and it does not require much computation or memory use. Although existing literature has analyzed its performance in various settings [1, 3, 8, 10], it is not known whether and when the policy is optimal. Our main result is the following theorem, which shows that Exact Scheduling minimizes the variance of service capacity under stationary job arrivals and strict demand and deadline constraints among distributed algorithms.

THEOREM 1. *Exact Scheduling (6) is the optimal solution of (5) and achieves the optimal value²*

$$\text{Var}(P) = \Lambda \mathbb{E} \left[\frac{\sigma^2}{\tau} \right].$$

Theorem 1 shows the achievable performance improvement by performing distributed service capacity control. If no control is applied, *i.e.*, $r_k(t) = \mathbf{1}\{t \in [a_k, a_k + \sigma_k]\}$, then $\mathbb{E}(P) = \text{Var}(P) = \Lambda \mathbb{E}[\sigma]$. By performing a distributed service capacity control, the stationary variance can be reduced by $\Lambda \mathbb{E}[\sigma(\tau - \sigma)/\tau]$, where $\tau - \sigma$ is the slack time (the amount of time left at job completion if a job is served at its maximum service rate upon arrival).

Given that we focus on distributed algorithms for scalability, it is important to understand how much performance degradation is incurred due to restricting ourselves to distributed policies compared to centralized policies. Centralized policies are the class of algorithms of the form

$$r_k(t) = w(k, t, A_t), \quad k \in \mathcal{V}, \quad (7)$$

where $A_t = \{(a_k, \sigma_k, \tau_k, x_k(t), y_k(t)) : a_k \leq t\}$ is the set that contains the information of jobs arriving before t , and w is a deterministic mapping from (k, t, A_t) to $r_k(t)$.

LEMMA 1. *Under any centralized policy of the form (7), the stationary variance of $P(t)$ is lower-bounded by*

$$\text{Var}(P) \geq \frac{\Lambda^2 \mathbb{E}[\sigma^2]^2}{4 \text{Var}(X)},$$

where $X(t)$ is the total amount of remaining service demand of jobs arriving before t .

Lemma 1 characterizes the trade-off between achieving a small variance of $X(t)$ and achieving a small variance of $P(t)$. An immediate consequence of Lemma 1 is a competitive-ratio-like bound of Exact Scheduling.

COROLLARY 1. *Let $\text{Var}(P)$ be the stationary variance of $P(t)$ that is attained by Exact Scheduling (6). Let $\text{Var}(P^*)$ be*

the minimum stationary variance attainable by any centralized algorithm (7) with the same $\text{Var}(X)$ as Exact Scheduling. Then, the following relation holds:

$$\frac{\text{Var}(P)}{\text{Var}(P^*)} \leq \frac{\mathbb{E}[\sigma^2/\tau] \mathbb{E}[\sigma^2 \tau]}{\mathbb{E}[\sigma^2]^2},$$

where all expectations on the right hand side are taken over the arrival distribution.

Corollary 1 suggests that if sojourn time τ is a deterministic random variable, then Exact Scheduling performs as well as the best centralized algorithm having the same $\text{Var}(X)$ as Exact Scheduling. One such example is when all jobs have identical sojourn times that equal their service demands, *i.e.*, $\tau_k = \tau_j$ and $\tau_k = \sigma_k$ for any $k, j \in \mathcal{V}$. In this system, due to constraints (1) and (2), both the optimal centralized policy and the optimal distributed policy is $r_k(t) = \mathbf{1}\{t \in [a_k, a_k + \tau_k]\}$.

4. REFERENCES

- [1] Giorgio C Buttazzo. *Hard real-time computing systems: predictable scheduling algorithms and applications*, volume 24. Springer Science & Business Media, 2011.
- [2] Mahdi Behrangrad. A review of demand side management business models in the electricity market. *Renewable and Sustainable Energy Reviews*, 47:270–283, 2015.
- [3] Andres Ferragut, Fernando Paganini, and Adam Wierman. Controlling the variability of capacity allocations using service deferrals. *ACM Transactions on Modeling and Performance Evaluation of Computing Systems (TOMPECS)*, 2(3):15, 2017.
- [4] Yorie Nakahira, Niangjun Chen, Lijun Chen, and Steven H Low. Smoothed least-laxity-first algorithm for ev charging. In *Proceedings of the Eighth International Conference on Future Energy Systems*, pages 242–251. ACM, 2017.
- [5] Nikhil Bansal, Tracy Kimbrel, and Kirk Pruhs. Speed scaling to manage energy and temperature. *Journal of the ACM (JACM)*, 54(1):3, 2007.
- [6] Shivendra S Panwar, Don Towsley, and Jack K Wolf. Optimal scheduling policies for a class of queues with customer deadlines to the beginning of service. *Journal of the ACM (JACM)*, 35(4):832–844, 1988.
- [7] John P Lehoczky. Using real-time queueing theory to control lateness in real-time systems. *ACM SIGMETRICS Performance Evaluation Review*, 25(1):158–168, 1997.
- [8] Chung Laung Liu and James W Layland. Scheduling algorithms for multiprogramming in a hard-real-time environment. *Journal of the ACM (JACM)*, 20(1):46–61, 1973.
- [9] François Baccelli and Bartłomiej Błaszczyszyn. *Stochastic Geometry and Wireless Networks, Volume I - Theory*. Now Publishers, 2009.
- [10] John Lehoczky, Lui Sha, and Ye Ding. The rate monotonic scheduling algorithm: Exact characterization and average case behavior. In *Real Time Systems Symposium, 1989., Proceedings.*, pages 166–171. IEEE, 1989.